

LOWER BOUND THEORY OF NONZERO ENTRIES IN SOLUTIONS OF ℓ_2 - ℓ_p MINIMIZATION

XIAOJUN CHEN*, FENGMIN XU[†], AND YINYU YE[‡]

Abstract. Recently, variable selection and sparse reconstruction are solved by finding an optimal solution of a minimization model where the objective function is the sum of a data-fitting term in ℓ_2 norm and a regularization term in ℓ_p norm ($0 < p < 1$). Since it is a nonconvex model, most algorithms for solving the problem can only provide an approximate local optimal solution, where nonzero entries in the solution cannot be identified theoretically. In this paper, we establish lower bounds for the absolute value of nonzero entries in every local optimal solution of the model, which can be used to identify zero entries precisely in any numerical solution. Therefore, we have developed a lower bound theorem to classify zero and nonzero entries in its every local solution. These lower bounds clearly show the relationship between the sparsity of the solution and the choice of the regularization parameter and norm, so that our theorem can be used for selecting desired model parameters and norms. Furthermore, we also develop error bounds for verifying accuracy of numerical solutions of the ℓ_2 - ℓ_p minimization model. To demonstrate applications of our theory, we propose a hybrid orthogonal matching pursuit-smoothing gradient (OMP-SG) method for solving the nonconvex, non-Lipschitz continuous ℓ_2 - ℓ_p minimization problem. Computational results show the effectiveness of the lower bounds for identifying nonzero entries in numerical solutions and the OMP-SG method for finding a high quality numerical solution.

Keywords: Variable selection, sparse solution, linear least-squares problem, ℓ_p regularization, smoothing approximation, first order condition, second order condition.

AMS Subject Classifications: 90C26, 90C46, 90C90

1. Introduction. We consider the following minimization problem

$$\min_{x \in R^n} \|Ax - b\|_2^2 + \lambda \|x\|_p^p, \quad (1.1)$$

where $A \in R^{m \times n}$, $b \in R^m$, $\lambda \in (0, \infty)$, $p \in (0, 1)$. Recently, minimization problem (1.1) attracted great attention in variable selection and sparse reconstruction [5, 7, 8, 9, 28]. The objective function of (1.1),

$$f(x) := \|Ax - b\|_2^2 + \lambda \|x\|_p^p$$

*Department of Applied Mathematics, The Hong Kong Polytechnic University, Kowloon, Hong Kong. Email: maxjchen@polyu.edu.hk. This author's work is supported partly by the Hong Kong Research Grants Council (PolyU5003/08P).

[†]School of Science and Department of Computer Science and Technology, Xi'an Jiaotong University, P.R. China. Email: fengminxu@mail.xjtu.edu.cn. This author's work is supported partly by the national 973 project of China 2007CB311002.

[‡]Department of Management Science and Engineering, School of Engineering, Stanford University, USA, and Visiting Professor of the Hong Kong Polytechnic University, Kowloon, Hong Kong. Email: yyeye@stanford.edu. This author is partially supported by DOE DE-SC0002009.

consists of a data fitting term $\|Ax - b\|_2^2$ and a regularization term $\lambda\|x\|_p^p$. Problem (1.1) is intermediate between the ℓ_2 - ℓ_0 minimization problem

$$\min_{x \in R^n} \|Ax - b\|_2^2 + \lambda\|x\|_0 \quad (1.2)$$

and the ℓ_2 - ℓ_1 minimization problem

$$\min_{x \in R^n} \|Ax - b\|_2^2 + \lambda\|x\|_1, \quad (1.3)$$

in the sense

$$\|x\|_0 = \sum_{\substack{i=1 \\ x_i \neq 0}}^n |x_i|^0, \quad \|x\|_p^p = \sum_{i=1}^n |x_i|^p, \quad \text{and} \quad \|x\|_1 = \sum_{i=1}^n |x_i|. \quad (1.4)$$

Naturally, one expects that using the ℓ_p norm¹ in the regularization term can find sparser solution than using the ℓ_1 norm, which was evidenced in extensive computational studies [5, 7, 8, 9, 28]. However, some major theoretical issues remain open. Is there any theoretical justification for solving minimization problem (1.1) with $p < 1$? What are the solution characteristics of (1.1)? Is there theory to dictate the choice of the regularization parameter λ and norm p ? Our first main contribution of this paper is to answer these questions. We establish lower bounds for the absolute value of nonzero entries in every local optimal solution of (1.1) only when $p < 1$. Therefore, we have developed a lower bound theorem to classify zero and nonzero entries in every local solution of (1.1). These lower bounds clearly show the relationship between the sparsity of the solution and the choice of the regularization parameter and norm, so that the theorem can be used to guide the selection of desired model parameters and norms in (1.1). It can be also used to identify zero and nonzero entries in the numerical optimal solution.

More specifically, using the second order necessary condition for a local minimizer, we present a component-wise lower bound

$$L_i = \left(\frac{\lambda p(1-p)}{2\|a_i\|^2} \right)^{\frac{1}{2-p}} \quad (1.5)$$

for each nonzero entry x_i^* of any local optimal solution x^* of (1.1), that is,

$$\text{for any } i \in \mathcal{N}, \quad L_i \leq |x_i^*|, \quad \text{if } x_i^* \neq 0,$$

which is equivalent to the following statement

$$\text{for any } i \in \mathcal{N}, \quad x_i^* \in (-L_i, L_i) \Rightarrow x_i^* = 0.$$

Here, $\mathcal{N} = \{1, \dots, n\}$ and a_i is the i th column of the matrix A . We show that the columns $\{a_i \mid i \in \text{support}(x^*)\}$ are linearly independent, which implies that $\|x^*\|_0 \leq m$ and the ℓ_2 - ℓ_p minimization problem (1.1) has a finite number of local minimizers.

Most minimization algorithms are descent-iterative in nature, that is, starting from an initial point x^0 they generate a sequence of points x^k , $k = 0, 1, \dots$, such

¹ $\|x\|_p$ ($0 < p < 1$) is a quasi-norm which satisfies the norm axioms except the triangle inequality. We call $\|x\|_p$ a norm for simplicity.

that the objective values $f(x^k)$ are strictly decreasing along the sequence. Thus, any local minimizer, including the global minimizer, that a descent algorithm may find must be in the level set $\{x : f(x) \leq f(x^0)\}$, and the set must contain at least one local minimizer. Therefore, in both theory and practice, one may be only interested in the minimizers satisfying $f(x) \leq f(x^0)$. Specifically, for our problem, the zero vector $x^0 = 0$ would be a trivial initial point for (1.1) with $f(0) = \|b\|^2$, the least squares solution of $\min_x \|Ax - b\|$ is another choice, and so is a point generated by any heuristic procedure such as the Orthogonal Matching Pursuit method. Based on this observation, we use the first necessary condition for a local minimizer to present a lower bound

$$L = \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}} \right)^{\frac{1}{1-p}} \quad (1.6)$$

for the absolute value of nonzero entries in a local optimal solution x^* of (1.1), which satisfies $f(x^*) \leq f(x^0)$, that is,

$$\text{for any } i \in \mathcal{N}, \quad L \leq |x_i^*|, \quad \text{if } x_i^* \neq 0.$$

Moreover, we show that the number of nonzero entries in every local optimal solution x^* satisfying $f(x^*) \leq f(x^0)$ is bounded by

$$\|x^*\|_0 \leq \min \left(m, \frac{f(x^0)}{\lambda L^p} \right).$$

The lower bounds in (1.5) and (1.6) are not only useful for identification of zero entries in local optimal solutions from approximation ones, but also for selection of the regularization parameter λ and norm $\|\cdot\|_p$. In particular, for a given norm $\|\cdot\|_p$, the lower bounds can help us to choose the regularization parameter λ for controlling the sparsity level of the solution. On the other hand, for a given λ , the lower bounds can also help us to understand the ℓ_2 - ℓ_p problem with different values $p \in (0, 1)$.

We need to mention that Nikolova [22] proved a similar bound based on the second order condition in a different context. Nikolova's result is important: it shows that using non-convex potential functions is good for piecewise constant image restoration. However, the result has not been used in practical algorithms, because one needs to solve an optimization problem to construct the bound. On the other hand, our first and second order bounds have explicit close forms and are easily checkable. Moreover, we have found that the bound based on the first order condition seems more effective in practice, especially when a good initial point x^0 is chosen.

Our second main contribution is on some numerical issues for solving (1.1). The ℓ_p norm $\|\cdot\|_p$ for $0 < p < 1$ is neither convex nor Lipschitz continuous. Solving the nonconvex, non-Lipschitz continuous minimization problem (1.1) is difficult.

Most optimization algorithms are only efficient for smooth and convex problems. Nevertheless, some algorithms for nonsmooth and nonconvex optimization problems have been developed recently [4, 12, 20, 29]. However, the Lipschitz continuity remains a necessary condition to define the Clarke subgradient in these algorithms. To overcome the non-Lipschitz continuity, some approximation methods have been considered for solving (1.1). For example, at the k th iteration, replacing $\|x\|_p^p$ by the

following terms [5, 7, 23]

$$\sum_{i=1}^n \frac{x_i^2}{((x_i^{k-1})^2 + \varepsilon_i)^{1-p/2}}, \quad \sum_{i=1}^n (|x_i| + \varepsilon_i)^p \quad \text{or} \quad \sum_{i=1}^n \frac{|x_i|}{(|x_i^{k-1}| + \varepsilon_i)^{1-p}}.$$

Here $\varepsilon \in R^n$ is a small positive vector. The question is: are there error bounds for verifying accuracy of numerical solutions of these approximation methods? We have resolved this question by developing several error bounds.

More precisely, we consider smoothing methods for nonconvex, nonsmooth optimization problems, for example, the smoothing gradient method [29]. We choose a smoothing function $s_\mu(t)$ of $|t|$, such that s_μ^p is continuously differentiable for any fixed scalar $\mu > 0$ and satisfies

$$0 \leq (s_\mu(t))^p - |t|^p \leq \left(\frac{\mu}{2}\right)^p.$$

See section 3. Let the smoothing objective function of f be

$$f_\mu(x) := \|Ax - b\|_2^2 + \sum_{i=1}^n (s_\mu(x_i))^p.$$

We can show that the solution x_μ^* of the smoothing nonconvex minimization problem

$$\min_{x \in R^n} f_\mu(x) \tag{1.7}$$

converges to a solution of (1.1) as $\mu \rightarrow 0$. For some small $\mu < \frac{L}{2}$, let

$$(\bar{x}_\mu^*)_i = \begin{cases} 0 & \text{if } |(x_\mu^*)_i| \leq \mu \\ (x_\mu^*)_i & \text{otherwise.} \end{cases}$$

We show that there is x^* in the solution set of (1.1) such that

$$(\bar{x}_\mu^*)_i = 0 \quad \text{if and only if} \quad x_i^* = 0, \quad i \in \mathcal{N}$$

and

$$\|\bar{x}_\mu^* - x^*\| \leq \kappa \|\nabla f_\mu(\bar{x}_\mu^*)\|, \tag{1.8}$$

where κ is a computable constant.

To demonstrate the significance of the absolute lower bounds (1.5), (1.6) and error bounds (1.8), we propose a hybrid orthogonal matching pursuit-smoothing gradient (OMP-SG) method for the nonconvex, non-Lipschitz ℓ_2 - ℓ_p minimization problem (1.1). We first use the orthogonal matching pursuit method to select candidates of nonzero entries in the solution. Next we use the smoothing gradient method in [29] to find an approximate solution of (1.1). Both before and after the SG method, we use the lower bound theory to identify zero entries in the solution.

Our preliminary numerical results show that using OMP-SG with elimination of small entries in the numerical solution by the lower bounds for ℓ_2 - ℓ_p minimization problem (1.1) can provide more sparse solutions with smaller predictor error compared with several well-known approaches for variable selection.

This paper is organized as follows. In section 2, we present absolute lower bounds (1.5) and (1.6) for nonzero entries in local solutions of ℓ_2 - ℓ_p minimization problem (1.1). In section 3, we present the computable error bound (1.8) for numerical solutions. In section 4, we give the hybrid OMP-SG method for solving the ℓ_2 - ℓ_p minimization problem (1.1). Numerical results are given to demonstrate the effectiveness of the lower bounds, the error bounds and the OMP-SG method.

Notations Throughout the paper, $\|\cdot\|$ denotes the ℓ_2 norm and $|\cdot|$ denotes the vector of the componentwise absolute value. For any $x, y \in R^n$, $x \cdot y$ represents the vector $(x_1y_1, \dots, x_ny_n)^T$ and $x^T y$ denotes the inner product. Let \mathcal{X}_p^* denote the set of local minimizers of (1.1). For a vector $x \in R^n$, $\text{support}(x) = \{i \in \mathcal{N} \mid x_i \neq 0\}$ denotes the support set of x .

2. Lower bounds for nonzero entries in solutions. In this section we present two lower bounds for nonzero entries in local solutions of ℓ_2 - ℓ_p minimization problem (1.1).

Since $f(x) \geq \lambda \|x\|_p^p$, the objective function $f(x)$ is bounded below and $f(x) \rightarrow \infty$ if $\|x\| \rightarrow \infty$. Moreover, the set \mathcal{X}_p^* of local minimizers of (1.1) is nonempty and bounded.

THEOREM 2.1. *(The second order bound) Let $L_i = \left(\frac{\lambda p(1-p)}{2\|a_i\|^2}\right)^{\frac{1}{2-p}}$, $i \in \mathcal{N}$. Then for any $x^* \in \mathcal{X}_p^*$, the following statements hold.*

(1)

$$\text{for any } i \in \mathcal{N}, \quad x_i^* \in (-L_i, L_i) \quad \Rightarrow \quad x_i^* = 0.$$

(2) *The columns of the sub-matrix $B := A_\Lambda \in R^{m \times |\Lambda|}$ of A are linearly independent, where $\Lambda = \text{support}(x^*)$, and $|\Lambda| = \|x^*\|_0$ is the cardinality of the set Λ .*

(3)

$$\|B^T A(x^* - b)\| \leq \frac{\lambda p}{2} \cdot \sqrt{\|x^*\|_0} \left(\min_{1 \leq i \leq \|x^*\|_0} L_i \right)^{p-1}.$$

In particular, If $\|a_i\| = 1$ for all $i \in \mathcal{N}$ (that is, A is column-wise normalized), then

$$\|B^T A(x^* - b)\| \leq \sqrt{\|x^*\|_0} \left(\frac{\lambda p}{2}\right)^{\frac{1}{2-p}} \left(\frac{1}{1-p}\right)^{\frac{1-p}{2-p}}.$$

(4)

$$\|x^*\| \leq \|(B^T B)^{-1} B^T b\| + \frac{\lambda p}{2} \|(B^T B)^{-1}\| \left(\min_{1 \leq i \leq |\Lambda|} L_i \right)^{p-1}.$$

If $\|a_i\| = 1$ for all $i \in \mathcal{N}$, then

$$\|x^*\| \leq \|(B^T B)^{-1} B^T b\| + \|(B^T B)^{-1}\| \left(\frac{\lambda p}{2}\right)^{\frac{1}{2-p}} \left(\frac{1}{1-p}\right)^{\frac{1-p}{2-p}}.$$

Proof. For $x^* \in \mathcal{X}_p^*$, with $\|x^*\|_0 = k$, without loss of generality, we assume

$$x^* = (x_1^*, \dots, x_k^*, 0, \dots, 0)^T.$$

Let $z^* = (x_1^*, \dots, x_k^*)^T$ and $B \in R^{m \times k}$ be the submatrix of A , whose columns are the first k columns of A . Define a function $g: R^k \rightarrow R$ by

$$g(z) = \|Bz - b\|^2 + \lambda \|z\|_p^p.$$

We have

$$f(x^*) = \|Ax^* - b\|^2 + \lambda \|x^*\|_p^p = \|Bz^* - b\|^2 + \lambda \|z^*\|_p^p = g(z^*).$$

Since $|z_i^*| > 0, i = 1, \dots, k$, g is continuously differentiable at z^* . Moreover, in a neighborhood of x^* ,

$$\begin{aligned} g(z^*) = f(x^*) &\leq \min\{f(x) \mid x_i = 0, i = k+1, \dots, n\} \\ &= \min\{g(z) \mid z \in R^k\}, \end{aligned}$$

which implies that z^* is a local minimizer of the function g . Hence the second order necessary condition for

$$\min_{z \in R^k} g(z) \tag{2.1}$$

holds at z^* .

(1) The second order necessary condition at z^* gives that the matrix

$$2B^T B + \lambda p(p-1)\text{diag}(|z^*|^{p-2})$$

is positive semi-definite. Therefore, we obtain

$$2e_i^T B^T B e_i + \lambda p(p-1)|z_i^*|^{p-2} \geq 0, \quad i = 1, \dots, k$$

where e_i is the i th column of the identity matrix of $R^{k \times k}$.

Note that $\|a_i\|^2 = e_i^T B^T B e_i$. We find that

$$|z_i^*|^{p-2} \leq \frac{2\|a_i\|^2}{\lambda p(1-p)}, \quad i = 1, \dots, k$$

which implies that

$$|z_i^*| \geq \left(\frac{\lambda p(1-p)}{2\|a_i\|^2} \right)^{\frac{1}{2-p}} = L_i, \quad i = 1, \dots, k.$$

Hence for any $x^* \in \mathcal{X}_p^*$, if $x_i^* \neq 0, i \in \mathcal{N}$, then $|x_i^*| \geq L_i$. This is equivalent to that if $x_i^* \in (-L_i, L_i), i \in \mathcal{N}$, then $x_i^* = 0$.

(2) Since the matrix $2B^T B + \lambda p(p-1)\text{diag}(|z^*|^{p-2})$ is positive semi-definite, and $\lambda p(p-1)\text{diag}(|z^*|^{p-2})$ is negative definite, the matrix $B^T B$ must be positive definite. Hence the columns of B must be linearly independent.

(3) Since z^* is a local minimizer of g , the first order necessary condition must hold at z^* . Hence, we find, with $Bz^* = Ax^*$,

$$\|B^T(Ax^* - b)\| = \|B^T(Bz^* - b)\| = \frac{\lambda p}{2} \|z^*\|^{p-1} \leq \frac{\lambda p}{2} \cdot \sqrt{|\Lambda|} \left(\min_{1 \leq i \leq |\Lambda|} L_i \right)^{p-1}.$$

If $\|a_i\| = 1$ for all $i \in \mathcal{N}$, then $L_i = \left(\frac{\lambda p(1-p)}{2} \right)^{\frac{1}{2-p}}$ for all $i \in \Lambda$, which implies (3).

(4) The first order necessary condition for (2.1) yields

$$2B^T Bz^* = 2B^T b - \lambda p |z^*|^{p-1} \text{sign}(z^*).$$

From (1) and (2) of this theorem, we know that $|z_i^*| \geq L_i$ and $B^T B$ is nonsingular. Hence, we obtain the desired results. For the case where $\|a_i\| = 1$, $i \in \mathcal{N}$, we have

$$\begin{aligned} \|x^*\| = \|z^*\| &\leq \|(B^T B)^{-1} B^T b\| + \|(B^T B)^{-1}\| \frac{1}{2} \lambda p \|z^*\|^{p-1} \\ &\leq \|(B^T B)^{-1} B^T b\| + \|(B^T B)^{-1}\| \left(\frac{\lambda p}{2} \right)^{\frac{1}{2-p}} \left(\frac{1}{1-p} \right)^{\frac{1-p}{2-p}}. \end{aligned}$$

□

COROLLARY 2.2. *The set \mathcal{X}_p^* of local minimizers of problem (1.1) has a finite number of elements. Moreover, we have*

$$\mathcal{X}_p^* \subseteq \left\{ x \mid \|x\| \leq \sigma \|A^T b\| + \sigma \frac{\lambda p}{2} \left(\min_{1 \leq i \leq |\Lambda|} L_i \right)^{p-1} \right\},$$

where

$$\sigma = \max\{ \|(B^T B)^{-1}\| \mid B \in R^{m \times k}, \text{rank}(B) = k, \text{ } B \text{ lies in the columns of } A \}$$

and $k = \text{rank}(A) \leq \min(m, n)$.

Proof. From (2) of Theorem 2.1, we find that \mathcal{X}_p^* has a finite number of elements as there are at most $\binom{n}{m}$ possible matrices B , and the linear independence of the columns guarantees that at most one local minimizer exists for each matrix.

It is known that for any two sets $\{\hat{a}_i, i = 1, \dots, \ell\} \subseteq \{a_i, i = 1, \dots, k\}$, the matrices $\hat{B} \in R^{m \times \ell}$ and $B \in R^{m \times k}$ whose columns lie on the two sets, respectively, satisfy

$$\lambda_{\min}(B^T B) \leq \lambda_{\min}(\hat{B}^T \hat{B}),$$

where λ_{\min} denotes the smallest eigenvalue. Thus $\|(B^T B)^{-1}\| \leq \sigma$ for any matrix B arising from an $x^* \in \mathcal{X}_p^*$.

From (4) of Theorem 2.1, and $\|B^T b\| \leq \|A^T b\|$, we find the closed ball containing \mathcal{X}_p^* in this corollary. □

Remark 2.1. Note that Theorem 2.1 holds for all local minimizers of (1.1). Result (1) of Theorem 2.1 presents a lower bound theory of nonzero entries in local minimizers of (1.1). Result (2) implies that columns of A corresponding to nonzero entries of x^*

must form a basis as long as $0 < p < 1$, while bound (3) shows that x^* approaches the least squares solution of $\min_x \|Ax - b\|$ (restricted to the support of x^*) as $\lambda \rightarrow 0$. Corollary 2.2 points out that (1.1) has a finite number of local minimizers, and presents a closed ball which contains all local minimizers of (1.1), and an upper bound for all nonzero entries in any local minimizer.

As we mentioned before, most minimization algorithms are descent-iterative in nature, that is, they generate a sequence of points x^k , $k = 0, 1, \dots$, such that the objective values $f(x^k)$ are strictly decreasing along the sequence. Thus, any local minimizer, including the global minimizer, that a descent algorithm may find must be in the level set $\{x : f(x) \leq f(x^0)\}$, where x^0 is any given initial point. Therefore, in both theory and practice, one may be only interested in the minimizers satisfying $f(x) \leq f(x^0)$. Indeed, our next theorem presents a lower bound theory of nonzero entries for any local minimizer x^* of (1.1) in $\{x : f(x) \leq f(x^0)\}$, and derives an upper bound on $\|x^*\|_0$. The upper bound indicates that for λ sufficiently large but finite, $\|x^*\|_0$ reduces to 0 for $0 < p < 1$, which means that $x^* = 0$ is the only global minimizer.

THEOREM 2.3. *(The first order bound) Let x^* be any local minimizer of (1.1) satisfying $f(x^*) \leq f(x^0)$ for an arbitrarily given initial point x^0 . Let $L = \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}} \right)^{\frac{1}{1-p}}$. Then we have*

$$\text{for any } i \in \mathcal{N}, \quad x_i^* \in (-L, L) \quad \Rightarrow \quad x_i^* = 0.$$

Moreover, the number of nonzero entries in x^* is bounded by

$$\|x^*\|_0 \leq \min \left(m, \frac{f(x^0)}{\lambda L^p} \right). \quad (2.2)$$

Proof. Suppose $f(x^*) \leq f(x^0)$, $x^* \in \mathcal{X}_p^*$. Then, we have

$$\begin{aligned} \|A^T(Ax^* - b)\|^2 &\leq \|A^T\|^2 \|Ax^* - b\|^2 \leq \|A^T\|^2 (\|Ax^* - b\|^2 + \lambda \|x^*\|_p^p) \\ &= \|A^T\|^2 f(x^*) \leq \|A^T\|^2 f(x^0). \end{aligned} \quad (2.3)$$

Recall the function g in the proof of Theorem 2.1. The first order necessary condition for

$$\min_{z \in \mathbb{R}^k} g(z)$$

at z^* gives

$$2B^T(Bz^* - b) + \lambda p(|z^*|^{p-1} \cdot \text{sign}(z^*)) = 0.$$

This, together with (2.3), implies

$$\lambda p \| |z^*|^{p-1} \| = 2 \|B^T(Bz^* - b)\| = 2 \|B^T(Ax^* - b)\| \leq 2 \|A^T(Ax^* - b)\| \leq 2 \|A\| \sqrt{f(x^0)}.$$

Therefore, we obtain

$$2 \|A\| \sqrt{f(x^0)} \geq \lambda p \| |z^*|^{p-1} \| \geq \lambda p \left(\min_{1 \leq i \leq k} |z_i^*| \right)^{p-1}.$$

Note that $p - 1 < 0$. We find

$$\min_{1 \leq i \leq k} |z_i^*| \geq \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}} \right)^{\frac{1}{1-p}} = L.$$

Hence, all nonzero components of x^* are no less than L . In other words, for $i \in \mathcal{N}$, if $x_i^* \in (-L, L)$ then $x_i^* = 0$.

Now we show the second part of the theorem. Again,

$$\lambda \|x^*\|_p^p \leq \|Ax^* - b\| + \lambda \|x^*\|_p^p = f(x^*) \leq f(x^0).$$

From the first part of this theorem, any nonzero entry of x^* is bounded from below by L . Thus, they together with (2) of Theorem 2.1, imply the desired bound in (2.2). \square

The lower bound in Theorem 2.1 depends on the parameters λ, p , and the matrix A , while the lower bound in Theorem 2.3 depends on λ, p, A and the initial objective value $f(x^0)$. In practice, we can take the maximum value of the two bounds to get a new bound. Moreover, in Theorem 2.3 one may simply set $x^0 = 0$, the trivial local minimizer of (1.1), (so that $f(x) \leq f(x^0) = \|b\|^2$), the minimizer of $\|Ax - b\|$, or a point generated by any heuristic procedure such as the Orthogonal Matching Pursuit method. It is worth noting that x^0 can be replaced by x^* in Theorem 2.3 and the theorem remains true.

The lower bound theory can be extended to the following problem

$$\min_{x \in R^n} \|Ax - b\|^2 + \lambda \sum_{i=1}^r \varphi(d_i^T x_i), \quad (2.4)$$

where $D \in R^{r \times n}$ is the first or second order difference matrix with rows d_i , and φ is a non-Lipschitz potential function; see Table 4.5. In fact, as we mentioned earlier, Nikolova [22] proved that there is $\theta > 0$ such that every local minimizer x^* of (2.4) satisfies

$$\text{either } |d_i^T x^*| = 0 \quad \text{or} \quad |d_i^T x^*| \geq \theta$$

by using the second order necessary condition for (2.4). However, the result has not been used in practical algorithms, because one needs to solve an optimization problem to construct θ . Nikolova [22] also stated that it is difficult to get an explicit solution from the optimization problem for constructing θ .

Lower bounds (1.5) and (1.6) clearly show the relationship between the sparsity of the solution and the choice of the regularization parameter λ and norm $\|\cdot\|_p$. Hence our lower bound theory can be used for selecting model parameters λ and p . In Figure 1, we show some properties of the function $L(\lambda, p) = (\lambda p(1-p))^{\frac{1}{2-p}}$ for $\lambda = (0, 10]$ and $p \in [0, 1]$.

From Figure 1, we can see clearly that for any given $\lambda > 0$, $(\lambda p(1-p))^{\frac{1}{2-p}}$ is a nonnegative and concave function of p on $[0, 1]$. It takes the minimum value at $p = 0$ and $p = 1$, for any $\lambda \in (0, 10]$.

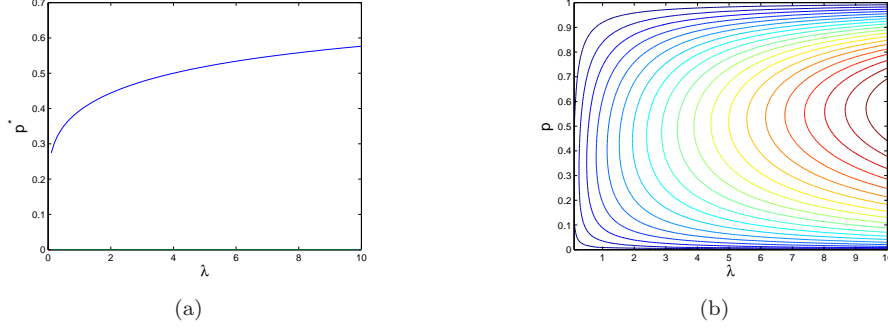


Fig. 2.1: (a) $p^*(\lambda) = \arg \max_{0 \leq p \leq 1} L(\lambda, p)$ (b) $L(\lambda, p) = (\lambda p(1-p))^{\frac{1}{2-p}}$

3. Error Bounds derived from lower bound theory. Smoothing approximations are widely used in optimization and scientific computing. In the following we consider a smoothing function of f and give a smooth version of Theorem 2.1 and Theorem 2.3.

For $\mu \in (0, \infty)$, let

$$s_\mu(t) = \begin{cases} |t| & |t| > \mu \\ \frac{t^2}{2\mu} + \frac{\mu}{2} & |t| \leq \mu. \end{cases}$$

Then $s_\mu(t)$ is continuously differentiable and

$$((s_\mu(t))^p)' = \begin{cases} p|t|^{p-1} \text{sign}(t) & |t| > \mu \\ p \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-1} \frac{t}{\mu} & |t| \leq \mu. \end{cases}$$

However, $s_\mu(t)$ is not twice differentiable at $t = \mu$. For $t \in (-\mu, \mu)$, the second derivative of $(s_\mu(t))^p$ satisfies

$$\begin{aligned} (s_\mu(t)^p)'' &= p(p-1) \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-2} \left(\frac{t}{\mu} \right)^2 + p \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-1} \frac{1}{\mu} \\ &\geq p(p-1) \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-2} \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right) \frac{1}{\mu} + p \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-1} \frac{1}{\mu} \\ &= p^2 \left(\frac{t^2}{2\mu} + \frac{\mu}{2} \right)^{p-1} \frac{1}{\mu} > 0. \end{aligned} \quad (3.1)$$

Hence $s_\mu^p(t)$ is strictly convex in $(-\mu, \mu)$. Moreover, from $s_\mu(t) = |t| \left(\frac{t^2 + \mu^2}{2\mu|t|} \right) \geq |t|$ and $0 = \arg \max_{t \in (-\mu, \mu)} (s_\mu(t) - |t|)$, we have that for any $t \in R$

$$0 \leq (s_\mu(t))^p - |t|^p \leq \left(\frac{\mu}{2} \right)^p. \quad (3.2)$$

Let

$$\psi_\mu(x) = (s_\mu(x_1), \dots, s_\mu(x_n))^T$$

and

$$\Psi_\mu(x) = \left(((s_\mu(x_1))^p)', \dots, ((s_\mu(x_n))^p)' \right)^T.$$

We define a smoothing approximation of the objective function $f(x)$

$$f_\mu(x) = \|Ax - b\|^2 + \lambda \|\psi_\mu(x)\|_p^p,$$

and consider the smooth minimization problem (1.7). The smoothing objective function f_μ is continuously differentiable in R^n , and strictly convex on the set $\{x \mid \|x\|_\infty \leq \mu\}$.

Let $\mathcal{X}_{p,\mu}^*$ denote the set of local minimizers of (1.7). By the definition of ψ_μ and (3.2) for any x we have

$$\lambda n \left(\frac{\mu}{2}\right)^p \geq f_\mu(x) - f(x) \geq 0.$$

Since $\|x\| \rightarrow \infty$ implies $f(x) \rightarrow \infty$, we deduce $f_\mu(x) \rightarrow \infty$ if $\|x\| \rightarrow \infty$. Moreover, for any $x \in R^n$, $\lim_{\mu \downarrow 0} f_\mu(x) = f(x)$. The following theorem presents the smooth version of the first and second order lower bounds.

THEOREM 3.1. *Let $L = \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}} \right)^{\frac{1}{1-p}}$ for an arbitrarily given initial point x^0 , and $L_i = \left(\frac{\lambda p(1-p)}{2\|a_i\|^2} \right)^{\frac{1}{2-p}}$, $i \in \mathcal{N}$.*

(1) *(The second order bound) For any $\mu > 0$ and any $x_\mu^* \in \mathcal{X}_{p,\mu}^*$, we have*

$$\text{for any } i \in \mathcal{N}, \quad (x_\mu^*)_i \in (-L_i, L_i) \quad \Rightarrow \quad |(x_\mu^*)_i| \leq \mu.$$

(2) *(The first order bound) For any $\mu > 0$ and any $x_\mu^* \in \mathcal{X}_{p,\mu}^*$ satisfying $f(x_\mu^*) \leq f(x^0)$, we have*

$$\text{for any } i \in \mathcal{N}, \quad (x_\mu^*)_i \in (-L, L) \quad \Rightarrow \quad |(x_\mu^*)_i| \leq \mu.$$

Proof. (1) Since $x_\mu^* \in \mathcal{X}_{p,\mu}^*$, the second order necessary condition for (1.7) implies that the matrix

$$\nabla^2 f_\mu(x_\mu^*) = 2A^T A + \lambda \Psi'_\mu(x)$$

is positive semi-definite. Suppose $|(x_\mu^*)_i| > \mu$ then from

$$e_i^T (2A^T A + \lambda \Psi'_\mu(x)) e_i = 2\|a_i\|^2 + \lambda p(p-1) |(x_\mu^*)_i|^{p-2} \geq 0,$$

we can get

$$|(x_\mu^*)_i| \geq \left(\frac{\lambda p(1-p)}{2\|a_i\|^2} \right)^{\frac{1}{2-p}} = L_i.$$

Since $\mu > 0$ and $x_\mu^* \in \mathcal{X}_{p,\mu}^*$ are arbitrarily chosen, we can claim that for any $\mu > 0$ and $x_\mu^* \in \mathcal{X}_{p,\mu}^*$, if $(x_\mu^*)_i \in (-L_i, L_i)$, $i \in \mathcal{N}$, then $|(x_\mu^*)_i| \leq \mu$.

(2) Since $x_\mu^* \in \mathcal{X}_{p,\mu}^*$, the first order necessary condition for (1.7) gives

$$\nabla f_\mu(x_\mu^*) = 2(A^T Ax_\mu^* - A^T b) + \lambda \Psi_\mu(x_\mu^*) = 0, \quad (3.3)$$

which, together with $f(x_\mu^*) \leq f(x^0)$, implies

$$\|\lambda \Psi_\mu(x_\mu^*)\|^2 \leq 4\|A^T\|^2(\|Ax_\mu^* - b\|^2 + \|x_\mu^*\|_p^p) = 4\|A\|^2 f(x_\mu^*) \leq 4\|A\|^2 f(x^0). \quad (3.4)$$

Suppose $|(x_\mu^*)_i| > \mu$ then

$$\lambda \|\Psi_\mu(x_\mu^*)\| \geq \lambda |\Psi_\mu(x_\mu^*)_i| = \lambda p |(x_\mu^*)_i|^{p-1}. \quad (3.5)$$

From (3.4) and (3.5) we can get

$$|(x_\mu^*)_i|^{p-1} \leq \frac{2\|A\|\sqrt{f(x^0)}}{\lambda p}.$$

Note that $p - 1 < 0$, we find

$$|(x_\mu^*)_i| \geq \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}} \right)^{\frac{1}{1-p}} = L.$$

Hence we can claim that for $i \in \mathcal{N}$ if $(x_\mu^*)_i \in (-L, L)$ then $|(x_\mu^*)_i| \leq \mu$. \square

The function f is not Lipschitz continuous. We define the first order necessary condition and the second order necessary condition for (1.1) as follows.

DEFINITION 3.2. For $x \in R^n$, let $X = \text{diag}(x)$.

(1) x is said to satisfy the first order necessary condition of (1.1) if

$$2XA^T(Ax - b) + \lambda p|x|^p = 0. \quad (3.6)$$

(2) x is said to satisfy the second order necessary condition of (1.1) if

$$2XA^TAX + \lambda p(p-1)\text{diag}(|x|^p) \quad (3.7)$$

is positive semi-definite.

Obviously, the zero vector in R^n satisfies the first and second necessary condition of (1.1).

Let $\{x_{\mu_k}\}$ denote a sequence with $\mu_k > 0$, $k = 1, 2, \dots$, and $\mu_k \rightarrow 0$ as $k \rightarrow \infty$.

THEOREM 3.3.

- (1) Let $\{x_{\mu_k}\}$ be a sequence of vectors satisfying the first order necessary condition of (1.7). Then any accumulation point of $\{x_{\mu_k}\}$ satisfies the first order necessary condition of (1.1).
- (2) Let $\{x_{\mu_k}\}$ be a sequence of vectors satisfying the second order necessary condition of (1.7). Then any accumulation point of $\{x_{\mu_k}\}$ satisfies the second order necessary condition of (1.1).
- (3) Let $\{x_{\mu_k}\}$ be a sequence of vectors being global minimizers of (1.7). Then any accumulation point of $\{x_{\mu_k}\}$ is a global minimizer of (1.1).

Proof. Let \bar{x} be an accumulation point of $\{x_{\mu_k}\}$. By working on a subsequence, we may assume that $\{x_{\mu_k}\}$ converges to \bar{x} . Let $X_{\mu_k} = \text{diag}(x_{\mu_k})$ and $\bar{X} = \text{diag}(\bar{x})$.

(1) From the first order necessary condition (3.3) of (1.7), we have

$$X_{\mu_k} \nabla f_{\mu_k}(x_{\mu_k}) = 2X_{\mu_k}(A^T A x_{\mu_k} - A^T b) + \lambda X_{\mu_k} \Psi_{\mu_k}(x_{\mu_k}) = 0.$$

By the definition of Ψ_{μ} , we have

$$(X_{\mu_k} \Psi_{\mu_k}(x_{\mu_k}))_i = p|x_{\mu_k}|_i^p, \quad \text{if } |x_{\mu_k}|_i > \mu_k$$

and

$$0 \leq (X_{\mu_k} \Psi_{\mu_k}(x_{\mu_k}))_i = p\left(\frac{(x_{\mu_k})_i^2}{2\mu_k} + \frac{\mu_k}{2}\right)^{p-1} \frac{(x_{\mu_k})_i^2}{\mu_k} \leq p\left(\frac{(x_{\mu_k})_i^2}{\mu_k}\right)^p \leq p|x_{\mu_k}|_i^p, \quad \text{if } |x_{\mu_k}|_i \leq \mu_k.$$

If $|x_{\mu_k}|_i \leq \mu_k$ for arbitrarily large k , then $\bar{x}_i = \lim_{k \rightarrow \infty} (x_{\mu_k})_i = 0$, and $\lim_{k \rightarrow \infty} (X_{\mu_k} \Psi_{\mu_k}(x_{\mu_k}))_i = 0$.

Therefore, we have

$$0 = 2 \lim_{k \rightarrow \infty} X_{\mu_k}(A^T A x_{\mu_k} - A^T b) + \lambda \lim_{k \rightarrow \infty} X_{\mu_k} \Psi_{\mu_k}(x_{\mu_k}) = \bar{X}(A^T A \bar{x} - A^T b) + \lambda p|\bar{x}|^p.$$

Hence \bar{x} satisfies the first order necessary condition of (1.1).

(2) From the second order necessary condition of (1.7), we have

$$X_{\mu_k} \nabla^2 f(x_{\mu_k}) X_{\mu_k} = 2X_{\mu_k} A^T A X_{\mu_k} + \lambda X_{\mu_k} \Psi'_{\mu_k}(x_{\mu_k}) X_{\mu_k}$$

is positive semi-definite. Using the definition of Ψ_{μ} and (3.1), we have

$$(X_{\mu_k} \Psi'_{\mu_k}(x_{\mu_k}) X_{\mu_k})_{ii} = p(p-1)|x_{\mu_k}|_i^p, \quad \text{if } |x_{\mu_k}|_i > \mu_k$$

and

$$\begin{aligned} 0 < (X_{\mu_k} \Psi'_{\mu_k}(x_{\mu_k}) X_{\mu_k})_{ii} &= p(p-1)\left(\frac{(x_{\mu_k})_i^2}{2\mu_k} + \frac{\mu_k}{2}\right)^{p-2} \frac{(x_{\mu_k})_i^4}{\mu_k^2} + p\left(\frac{(x_{\mu_k})_i^2}{2\mu_k} + \frac{\mu_k}{2}\right)^{p-1} \frac{(x_{\mu_k})_i^2}{\mu_k} \\ &\leq p\left(\frac{(x_{\mu_k})_i^2}{2\mu_k} + \frac{\mu_k}{2}\right)^{p-1} \frac{(x_{\mu_k})_i^2}{\mu_k} \leq p\left(\frac{(x_{\mu_k})_i^2}{2\mu_k} + \frac{(x_{\mu_k})_i^2}{2\mu_k}\right)^{p-1} \frac{(x_{\mu_k})_i^2}{\mu_k} \leq p\left(\frac{(x_{\mu_k})_i^2}{\mu_k}\right)^p \leq p\mu_k^p, \\ &\quad \text{if } |x_{\mu_k}|_i \leq \mu_k. \end{aligned}$$

Therefore, for any $y \in R^n$, we have

$$\begin{aligned} 0 &\leq \lim_{k \rightarrow \infty} y^T (2X_{\mu_k} A^T A X_{\mu_k} + \lambda X_{\mu_k} \Psi'_{\mu_k}(x_{\mu_k}) X_{\mu_k}) y \\ &= y^T \left(2 \lim_{k \rightarrow \infty} X_{\mu_k} A^T A X_{\mu_k} + \lambda \lim_{k \rightarrow \infty} X_{\mu_k} \Psi'_{\mu_k}(x_{\mu_k}) X_{\mu_k} \right) y \\ &= y^T (2\bar{X} A^T A \bar{X} + \lambda p(p-1) \text{diag}(|\bar{x}|^p)) y. \end{aligned}$$

Hence \bar{x} satisfies the second order necessary condition for (1.1).

(3) Let x^* be a global minimizer of (1.1). Then from the following three inequalities

$$f(x_{\mu_k}) \leq f_{\mu_k}(x_{\mu_k}) \leq f_{\mu_k}(x^*) \leq f(x^*) + \lambda n \left(\frac{\mu_k}{2}\right)^p,$$

we deduce that \bar{x} is a global minimizer of (1.1). \square

In the following, we present a computable error bound for KKT solutions (satisfying the first order necessary condition) of the smooth minimization problem (1.7) to approximate a KKT solution of the non-Lipschitz optimization problem (1.1).

Let $\mathcal{X}_{p,\mu}$ be the set of KKT solutions of (1.7) and \mathcal{X}_p be the set of KKT solutions of (1.1).

THEOREM 3.4. *Let $\{x_{\mu_k}\}$ be a sequence of vectors satisfying the first order necessary condition of (1.7) and $f(x_{\mu_k}) \leq f(x^0)$ for an arbitrarily given initial point x^0 . Then there is a $K > 0$, such that for any $k \geq K$, there is $x^* \in \mathcal{X}_p$ such that*

$$\Gamma_{\mu_k} := \{i \in \mathcal{N} \mid |(x_{\mu_k})_i| \leq \mu_k\} = \{i \in \mathcal{N} \mid |x_i^*| = 0\} =: \Gamma. \quad (3.8)$$

Define

$$(\bar{x}_{\mu_k}^*)_i = \begin{cases} 0 & i \in \Gamma \\ (x_{\mu_k})_i & i \in \mathcal{N} \setminus \Gamma. \end{cases} \quad (3.9)$$

Let B be the submatrix of A whose columns are indicated by $\mathcal{N} \setminus \Gamma$. Suppose $\lambda_{\min}(B^T B) > \frac{\lambda p(1-p)}{2} L^{p-2}$, then

$$\|\bar{x}_{\mu_k}^* - x^*\| \leq \|G^{-1}\| \|\nabla f_{\mu_k}(\bar{x}_{\mu_k}^*)\|, \quad (3.10)$$

where $G = 2B^T B + \lambda p(p-1)L^{p-2}I$, and $\lambda_{\min}(B^T B)$ denotes the smallest eigenvalue of the matrix $B^T B$.

Proof. Since the level set $\{x \mid f(x) \leq f(x^0)\}$ is bounded, the sequence $\{x_{\mu_k}\}$ is bounded. From (1) of Theorem 3.2, any accumulation point of $\{x_{\mu_k}\}$ is in \mathcal{X}_p . Hence we have $\lim_{k \rightarrow \infty} \text{dist}(x_{\mu_k}, \mathcal{X}_p) = 0$. This implies that there is $x^* \in \mathcal{X}_p$ such that $\lim_{k \rightarrow \infty} x_{\mu_k} = x^*$ and there is $K > 0$ such that for $k \geq K$, $\mu_k < \frac{L}{2}$,

$$\text{dist}(x_{\mu_k}, \mathcal{X}_p) = \|x_{\mu_k} - x^*\| < \frac{L}{2},$$

and $f(x^*) \leq f(x^0)$ hold. Then

$$|x_i^*| - |(x_{\mu_k})_i| \leq |x_i^* - (x_{\mu_k})_i| \leq \|x^* - x_{\mu_k}\| < \frac{L}{2}.$$

If $i \in \Gamma_{\mu_k}$, that is, $|(x_{\mu_k})_i| \leq \mu_k$, then we have

$$|x_i^*| < |(x_{\mu_k})_i| + \frac{L}{2} < L.$$

Assume that $x_i^* \neq 0$. From (3.6), we derive

$$\lambda p L^{p-1} < \lambda p |x_i^*|^{p-1} = 2|A^T(Ax^* - b)|_i \leq 2\|A^T(Ax^* - b)\| \leq 2\|A\|\sqrt{f(x^0)}$$

which implies $L > \left(\frac{\lambda p}{2\|A\|\sqrt{f(x^0)}}\right)^{1-p} = L$. This is a contradiction. Hence $|x_i^*| = 0$, that is, $i \in \Gamma$. We obtain that $\Gamma_{\mu_k} \subset \Gamma$.

On the other hand, if $i \in \Gamma$ then $x_i^* = 0$. We have

$$|(x_{\mu_k})_i| = |(x^* - x_{\mu_k})_i| \leq \|x^* - x_{\mu_k}\| < \frac{L}{2} < L.$$

From Theorem 3.1 we know $|(x_{\mu_k})_i| \leq \mu_k$, and thus $i \in \Gamma_{\mu_k}$. Hence $\Gamma \subset \Gamma_{\mu}$. We obtain (3.8).

Without loss of generality, we assume that $\mathcal{N} \setminus \Gamma = \{1, 2, \dots, r\}$. Define the function $g : R^r \rightarrow R$ by

$$g(z) = \|Bz - b\|_2^2 + \lambda \|z\|_p^p.$$

The first order necessary condition (3.6) at x^* yields

$$\nabla g(z^*) = 2B^T(Bz^* - b) + \lambda p |z^*|^{p-1} \cdot \text{sign}(z^*) = 0$$

at $z^* = (x_1^*, \dots, x_r^*)^T$. Furthermore, let $z_{\mu_k} = ((x_{\mu_k})_1, \dots, (x_{\mu_k})_r)^T$, then

$$\begin{aligned} \nabla g(z_{\mu_k}) &= \nabla g(z_{\mu_k}) - \nabla g(z^*) \\ &= 2B^T B(z_{\mu_k} - z^*) + \lambda p |z_{\mu_k}|^{p-1} \cdot \text{sign}(z_{\mu_k}) - \lambda p |z^*|^{p-1} \cdot \text{sign}(z^*). \end{aligned}$$

Note that $\text{sign}(z_{\mu_k}) = \text{sign}(z^*)$. By using the mean value theorem, we have

$$\begin{aligned} \nabla g(z_{\mu_k}) &= 2B^T B(z_{\mu_k} - z^*) + \lambda p \text{sign}(z_{\mu_k}) \cdot (|z_{\mu_k}|^{p-1} - |z^*|^{p-1}) \\ &= (2B^T B + \lambda p(p-1) D)(z_{\mu_k} - z^*), \end{aligned} \quad (3.11)$$

where $D \in R^{r \times r}$ is a diagonal matrix whose diagonal elements are $|\tilde{z}_{\mu_k}|_i^{p-2}$, where $(\tilde{z}_{\mu_k})_i$ is between $(z_{\mu_k})_i$ and z_i^* , $i = 1, 2, \dots, r$. Since $|(z_{\mu_k})_i| \geq L$, $|z_i^*| \geq L$, and $\text{sign}((z_{\mu_k})_i) = \text{sign}(z_i^*)$, we have $|\tilde{z}_{\mu_k}|_i \geq L$, $i = 1, 2, \dots, r$.

Since the matrix $2B^T B + \lambda p(p-1)D$ is symmetric, $0 < p < 1$ and $|\tilde{z}_{\mu_k}|_i \geq L$ for all $i \in \mathcal{N} \setminus \Gamma$, for any $z \in R^r$ with $\|z\| = 1$, we have

$$\begin{aligned} z^T (2B^T B + \lambda p(p-1)D)z &= z^T (2B^T B)z + \lambda p(p-1)z^T D z \\ &\geq 2z^T (B^T B)z + \lambda p(p-1)L^{p-2}\|z\|^2 \\ &\geq 2\lambda_{\min}(B^T B) + \lambda p(p-1)L^{p-2} \\ &> 0, \end{aligned}$$

where the last inequality uses the assumption of this theorem. Hence the matrix $2B^T B + \lambda p(p-1) D$ is invertible. We conclude from (3.9) and (3.11) that

$$\begin{aligned} \|\bar{x}_{\mu_k}^* - x^*\| &= \|z_{\mu_k} - z^*\| \leq \|(2B^T B + \lambda p(p-1) D)^{-1}\| \|\nabla g(z_{\mu_k})\| \\ &\leq \|(2B^T B + \lambda p(p-1) L^{p-2} I)^{-1}\| \|\nabla g(z_{\mu_k})\| \\ &= \|G^{-1}\| \|\nabla g(z_{\mu_k})\| \\ &\leq \|G^{-1}\| \|\nabla f_{\mu_k}(\bar{x}_{\mu_k}^*)\|, \end{aligned}$$

where the last inequality uses $\|\nabla g(z_{\mu_k})\| \leq \|\nabla f_{\mu_k}(\bar{x}_{\mu_k}^*)\|$, which can be shown as follows,

$$\begin{aligned}
\|\nabla g(z_{\mu_k})\| &= \|2B^T(Bz_{\mu_k}^* - b) + \lambda p|z_{\mu_k}|^{p-1} \cdot \text{sign}(z_{\mu_k})\| \\
&= \|2B^T(A\bar{x}_{\mu_k}^* - b) + \lambda p|z_{\mu_k}|^{p-1} \cdot \text{sign}(z_{\mu_k})\| \\
&= \|2B^T(A\bar{x}_{\mu_k}^* - b) + \lambda \Psi_{\mu_k}(z_{\mu_k})\| \\
&\leq \|2A^T(A\bar{x}_{\mu_k}^* - b) + \lambda \Psi_{\mu_k}(\bar{x}_{\mu_k}^*)\| \\
&= \|\nabla f_{\mu_k}(\bar{x}_{\mu_k}^*)\|,
\end{aligned}$$

where the inequality uses $(\bar{x}_{\mu_k}^*)_i = 0$ for $i \in \Gamma$ and $(s_{\mu_k}^p)'(0) = 0$. \square

4. Hybrid OMP-SG algorithm using lower bound theory. The lower bound theory can be applied to improve existing algorithms and develop new algorithms. To demonstrate the application, we use a hybrid Orthogonal Matching Pursuit-smoothing gradient (OMP-SG) method to solve the ℓ_2 - ℓ_p minimization problem (1.1). More specifically, we employ the OMP method to generate an initial point x^0 and its support, develop an SG method to further reduce the objective value of (1.1), and finally apply our theoretical result to purify the numerical solution by deleting its entries with small values. Our limited computational experiment in this section does not intend to develop a new algorithm for sparse reconstruction, but to show how our theory could improve any existing algorithm to achieve a higher quality performance.

The OMP algorithm is well-known in the literature of signal processing. The following algorithm is a standard version of the OMP algorithm [3], but has a different stop criterion.

ALGORITHM 1. Orthogonal Matching Pursuit(OMP)

Parameters: Given the $m \times n$ matrix A , the vector $b \in R^m$ and the error threshold β_0 .

Initialization: Initialize $k = 0$, and set

- the initial solution $x^0 = 0$.
- the initial residual $r^0 = b - Ax^0 = b$.
- the initial solution support $\Lambda_0 = \emptyset$.

Main Iteration: Increment k by 1 and perform the following steps:

- Find the index j_k that solves the optimization problem

$$j_k = \arg \max \frac{\|(Ax^{k-1} - b)^T a_j\|_2^2}{\|a_j\|}$$
 for $j \in \mathcal{N} \setminus \Lambda_{k-1}$.

- Let $\Lambda_k = \Lambda_{k-1} \cup \{j_k\}$.
- Compute x^k , the minimizer of $\|Ax - b\|_2^2$ subject to $\text{support}(x) = \Lambda^k$.
- Calculate the new residual $r^k = Ax^k - b$.
- If $\|A^T r^k\| < \beta_0$, stop, and let $\Lambda = \Lambda_k$.

Output: A point $x_{omp} := x^k$, a set $\Lambda = \text{support}(x_{omp})$ and a matrix $B = A_\Lambda \in R^{m \times |\Lambda|}$.

The smoothing gradient method (SG) [29] is a simple method for Lipschitz continuous but nonsmooth nonconvex minimization problems.

ALGORITHM 2. **Smoothing Gradient(SG)**

Step 1. Choose constants $\sigma, \rho \in (0, 1)$, and an initial point x^0 . Set $k = 0$.

Step 2. Compute the step size ν_k by the Armijo line search, where $\nu_k = \max\{\rho^0, \rho^1, \dots\}$ and ρ^i satisfies

$$f_{\mu_k}(x^k - \rho^i g_k) \leq f_{\mu_k}(x^k) - \sigma \rho^i g_k^T g_k.$$

Set $x^{k+1} = x^k - \nu_k g_k$. Here $g_k = \nabla f_{\mu_k}(x^k)$.

Step 3. If $\|\nabla f_{\mu_k}(x^{k+1})\| \geq n\mu_k$, then set $\mu_{k+1} = \mu_k$; otherwise, choose $\mu_{k+1} = \sigma\mu_k$.

Now we present the hybrid OMP-SG algorithm for solving ℓ_2 - ℓ_p minimization problem (1.1) with the lower bound L defined in (1.6).

ALGORITHM 3. **Hybrid OMP-SG**

Step 1. Using the OMP algorithm to get x_{omp} , $\Lambda = \text{support}(x_{omp})$ and $B = A_\Lambda \in \mathbb{R}^{m \times |\Lambda|}$.

Step 2. Using the SG algorithm with an initial point $x^0 = x_{omp}$ to find

$$y^* = \arg \min g(y) := \|By - b\|_2^2 + \lambda \|y\|_p^p.$$

Step 3. Output a numerical solution x^* , where

$$x_j^* = \begin{cases} y_j^* & |y_j^*| \geq L \quad \text{and} \quad j \in \Lambda, \\ 0 & \text{otherwise.} \end{cases}$$

Remark 4.1 We know that from (2.3) at a solution x^* of (1.1), $f(x^*) \leq f(x^0)$ and $\|A^T(Ax^* - b)\| \leq \|A\| \sqrt{f(x^0)}$ for any $x^0 \in \mathbb{R}^n$. From Theorem 2.1 and Theorem 2.3, the number of nonzero entries of x^* is less than $\kappa = \min\left(m, \frac{f(x^0)}{\lambda L^p}\right)$, and each nonzero entry satisfies $|x_i^*| \geq L$. In the hybrid OMP-SG method, we first choose candidates of columns of A which correspond to nonzero entries in a solution of (1.1) and use these candidates of columns to build a submatrix B . Based on Theorem 2.1 and Theorem 2.3, the number of columns of B is chosen slightly large than κ and the error threshold β_0 is slightly large than $\|A\| \sqrt{f(x^0)}$. Next, we use the globally convergent smoothing gradient method to find an approximate minimizer of the reduced problem $\min g(y)$. According to Theorem 3.4, we set some entries of the approximate solution to zero if their absolute values are less than L . It is worth noting that the lower bound theory is algorithms independent. For instance, we can replace the SG by the smoothing conjugate gradient (SCG) method [12] in Step 2 of the hybrid OMP-SG method to accelerate the algorithm, and have a hybrid OMP-SCG method.

Now we report numerical results to compare the performance of the hybrid OMP-SG method and OMP-SCG method for solving (1.1) with several other approaches to find sparse solutions. Our preliminary computational results indicate that the variable elimination according to our theory makes a significant difference. The computational test was conducted on a Philips PC (2.36 GHz, 1.96GB of RAM) with using Matlab 7.4.

We consider the following four approaches.

- LASSO: Solve the ℓ_2 - ℓ_1 problem (1.3) by the least squares algorithm (Lars) proposed in [16].

- ConApp: Solve the $\ell_2\text{-}\ell_p$ problem (1.1) with $p = \frac{1}{2}$ by using the following $\ell_2\text{-}\ell_1$ convex approximation [5]

$$\min \|Ax - b\|_2^2 + \lambda \sum_{i=1}^n \frac{|x_i|}{\sqrt{|x_i^{k-1}| + \varepsilon}} \quad (4.1)$$

at the k th iteration, where $\varepsilon > 0$ is a parameter. We use the Lars to solve (4.1).

- OMP-SG: Solve the $\ell_2\text{-}\ell_p$ problem (1.1) by the hybrid OMP-SG.
- OMP-SCG: Solve the $\ell_2\text{-}\ell_p$ problem (1.1) by the hybrid OMP-SCG.

4.1. Variable selection. This example is artificially generated and was firstly used in Tibshirani [25] to test the effectiveness of Lasso. The true solution is $x^* = (3, 1.5, 0, 0, 2, 0, 0, 0)^T$. We simulated 100 data sets consisting of m observations from the model

$$Ax = b + \sigma\eta,$$

where η is a noise vector generated by the standard normal distribution. We select three cases to discuss the performance of the three approaches LASSO, ConApp and OMP-SCG. The first case is $m = 40$, $\sigma = 3$, the second case is $m = 40$, $\sigma = 1$ and the last case is $m = 60$, $\sigma = 1$. We used 80 of the 100 data sets to select the variables, then tested the performance on the remaining 20. The Mean Squared Errors (MSE) over the test set are summarized in Table 4.1. The average number of correctly identified zero coefficient (ANZ) and the average number of the coefficients erroneously set to zero ($NANZ$) over test set are also presented in Table 4.1. In our numerical experiment, we used $p = 0.5$ and $\lambda \approx 1.1$ in the $\ell_2\text{-}\ell_p$ problem (1.1). From Table 4.1, we observe that OMP-SCG performs the best, followed by LASSO and ConApp.

Table 4.1: Results for variable selection

m	σ	Approach	MSE	ANZ	$NANZ$
40	3	LASSO	0.4730	4.77	0.23
		ConApp	0.4688	4.83	0.17
		OMP-SCG	0.4755	4.88	0.12
40	1	LASSO	0.1595	4.77	0.23
		ConApp	0.1541	4.86	0.14
		OMP-SCG	0.1511	4.91	0.09
60	1	LASSO	0.3582	4.92	0.08
		ConApp	0.3503	4.93	0.07
		OMP-SCG	0.3464	4.95	0.05

4.2. Signal reconstruction. The signal reconstruction has been studied extensively in the past decades [6, 15]. According to Donoho [15], the signal reconstruction can be solved by the $\ell_2\text{-}\ell_1$ model (1.3). In this subsection we apply the $\ell_2\text{-}\ell_p$ model with $p = 0.5$ to solve signal reconstruction problems.

Consider a real-valued, finite-length signal $x \in R^n$. Suppose x is T -sparse, that is, only T of the signal coefficients are nonzero and the others are zero. We use the following Matlab code to generate the original signal, a matrix A and a vector b .

```
x_or = zeros(n,1); q = randperm(n); x_or(q(1:T)) = 2*randn(T,1);
```

$A = \text{randn}(m,n); A = \text{orth}(A)'; b = A*x_{or} ;$

Our aim is to obtain good reconstructions of x with less nonzero entries. We applied OMP-SCG, LASSO and ConApp to reconstruct the signal. The error between the reconstructed signal and the original one is computed by 2-norm.

In Table 4.2.1, we present numerical results of three sets of signal examples with different values of L and λ . The CPU time is given in second. From Table 4.2.1, we observe that the three approaches can reconstruct the original signal with $n = 512, T = 60, m = 184$, while OMP-SCG has the highest accuracy. Moreover, the LASSO can not reconstruct the original signal with $n = 512, T = 60, m = 182$, but OMP-SCG and ConApp can reconstruct the original signal, while OMP-SCG has small error. Furthermore, if the original signal has $n = 512, T = 130, m = 225$, LASSO and ConApp algorithms can not reconstruct this signal, but OMP-SCG can reconstruct this signal with error=0.41. OMP-SG gives similar results as OMP-SCG, but uses more time.

Table 4.2.1: Results for signal reconstruction without noisy

Problem	LASSO	ConApp	OMP-SCG			
	(Error,Time)	(Error,Time)	L	λ	Error	Time
$n = 512$ $T = 60$ $m = 184$	$(5.33 \times 10^{-4},$ 0.653)	$(1.29 \times 10^{-5},$ 6.82)	0.8	0.002	1.12×10^{-16}	1.02
$n = 512$ $T = 60$ $m = 182$	(38.64, 0.43)	$(2.41 \times 10^{-5},$ 7.84)	0.7	0.001	1.03×10^{-16}	1.34
$n = 512$ $T = 130$ $m = 225$	(122.25, 0.69)	(119.43, 19.99)	0.00001	0.00006	0.41	4.03

Now we add noisy signals to the problem

$$b = A*x_{or} - w ;$$

where $w = \sigma\eta$ is independent identically distributed Gaussian noise with zero mean and variance σ^2 . We measure the quality of a reconstructed signal \hat{x} using the mean-square error(MSE), defined as $E[\|\hat{x} - x_{or}\|^2]$. To compare the capability of algorithms in recovering the original signals under noisy circumstance, we use the oracle estimator, defined as

$$x_{oracle} = \sigma^2 \text{tr}(A_{\Lambda}^T A_{\Lambda})^{-1}$$

where $\Lambda = \text{support}(x_{or})$.

For each algorithm, we calculated the ratio of the MSE of a reconstructed signal generated from the algorithm and the MSE of the oracle estimator and listed the results as ‘‘Ratio’’ in Table 4.2.2. The closer the ratio is to 1, the more robust is the algorithm. From Table 4.2.2, we can see that the Ratio of OMP-SCG is always closer to 1 than LASSO and ConApp.

Table 4.2.2: Results for signal reconstruction with noisy ($n = 512, T = 130, \sigma = 0.1$)

m	Method(s)	MSE	Ratio	CPU	m	Method(s)	MSE	Ratio	CPU
330	LASSO	3.71	1.46	3.2541	310	LASSO	5.34	1.77	1.7519
	ConApp	3.58	1.41	63.01		ConApp	4.10	1.36	60.83
	OMP-SCG	3.42	1.34	5.23		OMP-SCG	4.05	1.34	22.45
	<i>Oracle</i>	<i>2.5434</i>				<i>Oracle</i>	<i>3.0180</i>		
300	LASSO	5.30	1.77	2.3011	275	LASSO	6.1	1.75	2.01
	ConApp	4.04	1.35	69.12		ConApp	5.05	1.45	78.15
	OMP-SCG	3.97	1.33	23.42		OMP-SCG	4.94	1.41	18.83
	<i>Oracle</i>	<i>2.9845</i>				<i>Oracle</i>	<i>3.4877</i>		

4.3. Prostate cancer. The data set in this subsection is downloaded from the UCI Standard database [1] for the study of prostate cancer. The data set consists of the medical records of 97 patients who were about to receive a radical prostatectomy. The predictors are eight clinical measures: lcavol, lweight, age, lbph, svi, lcp, gleason and pgg45. Detailed explanation can be found in the UCI Standard database. This is a variable selection problem with $A \in R^{97 \times 8}$. One of our main aims is to identify which predictors are most significant in predicting the response.

The prostate cancer data were divided into two parts: a training set with 67 observations and a test set with 30 observations. The prediction error is the mean squared errors over the test set. The numerical results of Ridge regression [19] and Best Subset [2] were derived from [18]. In this example, we also select $p = 0.5$ in the ℓ_2 - ℓ_p model (1.1).

From Table 4.3 we find that OMP-SG and OMP-SCG succeed in finding three main factors and have smaller prediction accuracy than ConApp and LASSO. This implies that OMP-SG and OMP-SCG can find more sparse solution with smaller prediction error than LASSO.

Table 4.3: Results for prostate cancer

Parameter	LASSO	Ridge	Best Subset	ConApp	OMP-SG	OMP-SCG
x_1 (lcavol)	0.545	0.389	0.740	0.6187	0.6436	0.6436
x_2 (lweight)	0.237	0.238	0.367	0.2362	0.2804	0.2804
x_3 (lage)	0	-0.029	0	0	0	0
x_4 (lbph)	0.098	0.159	0	0.1003	0	0
x_5 (svi)	0.165	0.217	0	0.1858	0.1856	0.1857
x_6 (lcp)	0	0.026	0	0	0	0
x_7 (gleason)	0	0.042	0	0	0	0
x_8 (pgg45)	0.059	0.123	0	0	0	0
Number of nonzreo	5	8	2	4	3	3
Prediction error	0.478	0.5395	0.5723	0.468	0.4418	0.4419

Now we apply Theorem 3.4 to compute the error bound $\|G^{-1}\|\|\nabla f_\mu(\bar{x}_\mu^*)\|$ of \bar{x}_μ^* to $x^* \in \mathcal{X}_p^*$, for a given $\mu > 0$. We set $\mu < 0.01$ and $p = 0.5$. The numerical results are listed in Table 4.4.

Table 4.4: Error bounds for $\|\bar{x}_\mu^* - x^*\|$

μ	L	λ	error bound
0.001	0.015	0.1304	1.5793×10^{-5}
0.0001	0.0119	0.1164	5.7310×10^{-6}
0.00001	0.0119	0.1164	5.5721×10^{-6}

It is worth noting that the lower bound theory, the error bounds and the hybrid OMP-SCG method can be extended to

$$\min_{x \in R^n} \|Ax - b\|_2^2 + \sum_{i=1}^n \varphi(x_i), \quad (4.2)$$

where $\varphi : R_+ \rightarrow R$ is a potential function, e.g. [23], which includes (1.1) as a special case. Table 4.5 lists some well-used potential functions (left) and their extensions (right).

Table 4.5: Potential functions (PFs) where $\alpha \in (0, 1)$ is a parameter

	Convex	Non Lipschitz
f_1	$\varphi(t) = t $	$\varphi(t) = t ^p$
	Non convex	Non Lipschitz
f_2	$\varphi(t) = t ^p$	$\varphi(t) = (t ^p)^\alpha$
f_3	$\varphi(t) = \frac{\alpha t }{1 + \alpha t }$	$\varphi(t) = \frac{\alpha t ^p}{1 + \alpha t ^p}$
f_4	$\varphi(t) = \log(\alpha t + 1)$	$\varphi(t) = \log(\alpha t ^p + 1)$

The numerical results with different potential functions and $\alpha = 0.1699$ are listed in Table 4.6. We observe that choosing $p \leq 0.5$ seems good for this example, since using $p \leq 0.5$ can find three main factors with smaller prediction error than $p > 0.5$.

Table 4.6: Comparisons of different p with different PFs

p	$(L, \text{Number of nonzero, Prediction error})$			
	f_1	f_2	f_3	f_4
0.9	(0.0001, 4, 0.4754)	(0.011, 4, 0.473)	(2.500, 4, 0.475)	(2.040, 4, 0.474)
0.8	(0.0015, 4, 0.4740)	(0.013, 4, 0.468)	(1.990, 4, 0.474)	(1.851, 4, 0.474)
0.7	(0.0050, 4, 0.4741)	(0.012, 4, 0.465)	(1.755, 4, 0.474)	(1.550, 4, 0.474)
0.6	(0.0084, 4, 0.4661)	(0.015, 3, 0.446)	(1.545, 4, 0.475)	(1.344, 4, 0.475)
0.5	(0.0119, 3, 0.4419)	(0.016, 3, 0.445)	(1.420, 3, 0.477)	(1.200, 3, 0.483)
0.4	(0.0148, 3, 0.4456)	(0.014, 3, 0.445)	(1.480, 3, 0.477)	(1.114, 3, 0.484)
0.3	(0.0176, 3, 0.4429)	(0.012, 3, 0.443)	(1.590, 3, 0.484)	(1.190, 3, 0.483)
0.2	(0.0196, 3, 0.4359)	(0.018, 3, 0.443)	(1.955, 3, 0.483)	(1.240, 3, 0.482)

5. Final remark. Using the first and second order necessary condition for a local minimizer, we establish lower bounds for nonzero entries in any local optimal solution of a minimization model where the objective function is the sum of a data-fitting term in ℓ_2 norm and a regularization term in ℓ_p norm ($0 < p < 1$). This establishes a theoretical justification by “zeroing” those entries in an approximate solution whose values are small enough, and explanation why the model generates more sparse solutions when the norm parameter $p < 1$.

Moreover, the lower bounds clearly show the relationship between the sparsity of the solution and the choice of the regularization parameter and norm. These provide a systematic mechanism for selecting the model parameters, such as regularization weight λ and norm p . Based on these results, we propose a hybrid orthogonal matching pursuit-smoothing gradient (OMP-SG) method for the nonconvex, non-Lipschitz continuous ℓ_2 - ℓ_p minimization problem. Numerical results show that using the OMP-SG method to solve the ℓ_2 - ℓ_p minimization problem (1.1) can provide more sparse solutions with smaller predictor error compared with several well-known approaches

for variable selection.

Acknowledgment. We would like to thank Prof. Misha Kilmer and two referees for their valuable comments and suggestions for improving the presentation of this paper. We also like to thank Rick Chartrand, the late Paul Tseng and Weijun Zhou for their helpful comments.

REFERENCES

- [1] C. Blake and C. Merz, Repository of machine learning databases [DB/OL], Irvine, CA:University of California, Department of Information and Computer Science, 1998.
- [2] L. Breiman, Better subset regression using the nonnegative garrote, *Technometrics*, 37 (1995), pp. 373-384.
- [3] A. M. Bruckstein, D. L. Donoho and M. Elad, From sparse solutions of systems of equations to sparse modeling of signals and images, *SIAM Review*, 51 (2009), pp. 34-81.
- [4] J. V. Burke, A. S. Lewis and M. L. Overton, A robust gradient sampling algorithm for nonsmooth, nonconvex optimization, *SIAM J. Optim.*, 15 (2005), pp. 751-779.
- [5] E.J. Candes, M.B. Wakin, and S.P. Boyd, Enhancing sparsity by reweighted ℓ_1 minimization, *J. Fourier Anal. Appl.*, 14 (2008), pp. 877-905.
- [6] E. Candes, J. Romberg and T. Tao, Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information, *IEEE Trans. Inf. Theory*, 52 (2006), pp. 489-509.
- [7] R. Chartrand, Exact reconstructions of sparse signals via nonconvex minimization, *IEEE Signal Process. Lett.*, 14 (2007), pp. 707-710.
- [8] R. Chartrand, Nonconvex regularization for shape preservation, *IEEE International Conference on Image Processing (ICIP)*, 2007.
- [9] R. Chartrand and W. Yin, Iteratively reweighted algorithms for compressive sensing, in *33rd International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2008
- [10] S. S. Chen, D. L. Donoho and M. A. Saunders, Atomic decomposition by basis pursuit, *SIAM J. Sci. Comput.*, 20 (1998), pp. 33-61.
- [11] X. Chen and Y. Ye, On homotopy-smoothing methods for box-constrained variational inequalities, *SIAM J. Control Optim.* 37 (1999), pp. 589-616.
- [12] X. Chen and W. Zhou, Smoothing nonlinear conjugate gradient method for image restoration using nonsmooth nonconvex minimization, to appear in *SIAM J. Imaging Sciences*.
- [13] F. H. Clarke, Optimization and nonsmooth analysis, John Wiley & Sons, New York, (reprinted by SIAM, Philadelphia), 1990.
- [14] G. Davis, S. Mallat and M. Avellaneda, Adaptive greedy approximations, *J. Constructive Approx.*, 13 (1997), pp. 57-98.
- [15] D. L. Donoho, Compressed sensing, *IEEE Trans Information Theory*, 52 (2006), pp. 1289-1306.
- [16] B. Efron, T. Hastie, I. Johnstone and R. Tibshirani, Least angle regression, *Annals Statistics*, 23 (2004), pp. 407-499.
- [17] J. Fan and R. Li, Variable selection via nonconcave penalized likelihood and its oracle properties, *J. Amer. Statist. Assoc.*, 96 (2001), pp. 1348-1360.
- [18] T. Hastie, R. Tibshirani and J. Friedman, *The elements of statistical learning: data mining, Inference and Prediction*, New York: Springer Verlag, 2001.
- [19] A. E. Hoerl and R. Kennard, Ridge regression: biased estimation for nonorthogonal problems, *Technometrics*, 12 (1970), pp. 55-67.
- [20] K. C. Kiwiel, Convergence of the gradient sampling algorithm for nonsmooth nonconvex optimization, *SIAM J. Optim.*, 18 (2007), pp. 379-388.
- [21] S. Mallat and Z. Zhang, Matching pursuits with time-frequency dictionaries, *IEEE Trans. Signal Process.*, 41 (1993), pp. 3397-3415.
- [22] M. Nikolova, Analysis of the recovery of edges in images and signals by minimizing nonconvex regularized least-squares, *SIAM J. Multiscale Model. Simul.*, 4 (2005), pp. 960-991
- [23] M. Nikolova, M. K. Ng, S. Zhang and W. Ching, Efficient reconstruction of piecewise constant images using nonsmooth nonconvex minimization, *SIAM J. Imaging Sciences*, 1 (2008), pp. 2-25.
- [24] Y. C. Pati, R. Rezaifar and P. S. Krishnaprasad, Orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition, in *Proceedings of the 27th Asilomar Conference on Signals, Systems and Computers*, 1 (1993), pp. 40-44.

- [25] R. Tibshirani, Regression shrinkage and selection via the Lasso. *J Royal Statist Soc B*, 58 (1996), pp. 267-288.
- [26] J. A. Tropp and A. C. Gilbert, Signal recovery from random measurements via orthogonal matching pursuit, *IEEE Trans. Information Theory*, 53 (2007), pp. 4655-4667.
- [27] J. A. Tropp, A. C. Gilbert, S. Muthukrishnan and M. J. Strauss, Improved sparse approximation over quasi-incoherent dictionaries, in *Proceedings of the IEEE International Conference on Image Processing, Barcelona, 2003*, pp. 37-40.
- [28] Z. Xu, H. Zhang, Y. Wang and X. Chang, $L_{\frac{1}{2}}$ regularizer, *Science in China Series F-Inf Sci.*, 52 (2009), pp. 1-9.
- [29] C. Zhang and X. Chen, Smoothing projected gradient method and its application to stochastic linear complementarity problems, *SIAM J. Optim.*, 20 (2009), pp. 627-649.